

Porting and tuning R5F code on Leonardo

M. YAGI

National Institutes for Quantum Science and Technology (QST) Rokkasho Institute for Fusion Energy

5-field Reduced Drift MHD Model

Model consists of Vorticity equation, Generalized Ohm's law, Continuity equation, Parallel momentum equation, Electron temperature evolution equation: (F, A, n, v, T_e)

It is an extension of 4-field reduced drift MHD model by Hazeltine and Meiss, Phys. Rep. 121 (1985)1-164.

$$\begin{split} &\frac{d}{dt}\nabla_{\perp}^{2}F - \delta[\nabla_{\perp}p_{i};\nabla_{\perp}F] = -\nabla_{\parallel}\nabla_{\perp}^{2}A - [\Omega,p] + \mu_{\perp}^{i}\nabla_{\perp}^{4}F \\ &\frac{d}{dt}\left(A - \delta^{2}\frac{m_{e}}{m_{i}}\nabla_{\perp}^{2}A\right) = -\nabla_{\parallel}(\phi - \delta p_{e}) + \eta_{\parallel}\nabla_{\perp}^{2}A - 4\mu_{\perp}^{e}\delta^{2}\frac{m_{e}}{m_{i}}\nabla_{\perp}^{4}A \\ &\frac{d}{dt}n + \beta\frac{d}{dt}p = \beta[\Omega,\phi - \delta p_{e}] - \beta\nabla_{\parallel}(v + \delta\nabla_{\perp}^{2}A) + \eta_{\perp}\beta\nabla_{\perp}^{2}p \\ &\frac{d}{dt}v = -\nabla_{\parallel}p + 4\mu_{\perp}^{i}\nabla_{\perp}^{2}v \\ &\frac{3}{2}\frac{dT_{e}}{dt} - \frac{\beta_{e}}{\beta}\frac{dn}{dt} = -\alpha_{T}\delta\beta_{e}\nabla_{\parallel}\nabla_{\perp}^{2}A + \varepsilon^{2}\chi_{e\parallel}\nabla_{\parallel}^{2}T_{e} + \chi_{e\perp}\nabla_{\perp}^{2}T_{e} \\ &\frac{d}{dt} = \frac{\partial}{\partial t} + [\phi,], \nabla_{\parallel} = \nabla_{\parallel}^{(0)} - [A,], F = \phi + \delta p_{i}, \delta p_{i} = \delta_{i}n, \delta p_{e} = \delta T_{e} + \delta_{e}n, \delta_{i} = \frac{\beta_{i}}{\beta}\delta, \delta_{e} = \frac{\beta_{e}}{\beta}\delta \end{split}$$



Application of 5-field Reduced MHD Model

Control of particle transport is an important issue for burn-control in ITER and DEMO.



S. Matsuda, NIFS-MEMO-80 (2017)



• Hollow density profile is often seen after pellet injection or gas-puff.

 \Rightarrow inversed density gradient appears in the edge region.

Particle pinch mechanism for inverted density gradient in semi-collisional regime [Miki APTWG ' 17]

ITG/TEM stability is investigated for inverted density gradient in collisionless limit [Du '17]

Parallel Numerical Algorithms (1D Domain Decomposition)





Cost Evaluation on HPE SGI8600 (NVIDIA V100)





Case with simply replacing FFTW with cuFFT

FFTW

```
call dfftw_plan_dft_c2r_2d(plan(1),NY,NZ,...)
call dfftw_plan_dft_c2r_2d(plan(2),NY,NZ,...)
    :
DO I=S,E
    DO L=1,N
      call dfftw_execute(plan(L))
    END DO
END DO
```



PCIバス (15GB/sec程度)

cuFFT

```
istat = cufftPlan2d(plan(1),NZ,NY,...)
istat = cufftPlan2d(plan(2),NZ,NY,...)
:
DO I=S,N
W1C_d=W1C
W2C_d=W2C
:
istat = cufftExecZ2D(plan(1),W1C_d,W1R_d)
istat = cufftExecZ2D(plan(2),W2C_d,W2R_d)
:
W1R=W1R_d
W2R=W2R_d
END DO
```

1 node 40 process 1 thread 1GPU





It is found that elapse time for data transfer between CPU and GPU is dominant in FFT area.

FFT library of GPU is faster than that of CPU: cuFFT 0.22sec vs FFTW 9.64 sec.



!\$acc enter data
create(W1C(1:LYMED,1:KZ
MWD,S:E),...)

!\$acc update
device(DXVOLK_T(S:E,1:ND
M),...)

!\$acc host_data use_device(W1C(1:LYMWD,1:K ZMWD,S:E),...)



!\$acc exit data delete(W1C,...)

!\$acc data present(& \$acc W1C(1:LYMWD,1:KZMWD,S:E), W2C(1:LYMWD,1:KZMWD,S:E),W3C(1:LYMWD,1:KZMWD,S:E),W4C(1:LYMWD,1:KZMWD,S:E), DO L=S1.E1 DO I=1.IRMAXM VOLRHS(I,L)=0.0D0 END DO END DO IF(IST ==0) THEN istat = cufftPlan2d(cuplan(1), KZMWD, KYMWD, CUFFT Z2D) IST=1 if (myid.eq.0) write(*,*) 'time fft plan:', MPI WTIME() - t0 ! okada ENDIF **!\$acc kernels** !\$acc loop collapse(3) DO I=S.E DO N=1.KZMWD DO M=1,LYMWD $W1C(M, \dot{N}, I) = (0.0D0, 0.0D0)$ END DO END DO END DO !\$acc loop independent private(IY,IZ) collapse(2) DO I=S.E **Execution on GPU** DO L=1,NDM IY=IPY(L) IZ=IPZ(L) $W1C(I\dot{Y},I\dot{Z}) = DXVOLK_T(I,L)$ END DO END DO Sacc end kernels Modified subroutine TMRHS

```
!$acc host data
use_device(W1C(1:LYMWD,1:KZMWD,S:E),W2C(1:LYMWD,1:KZMWD,S:E),W3C(1:LYMWD,1:KZMWD,S:E),W4C(1:LYMWD,1:KZMWD,S:E),&
DO I=S F
 istat = cufftExecZ2D(cuplan(1), W1C(1,1,I), W1R(1,1,I))
END DO
!$acc end host data
!$acc kernels
!Sacc loop
DO I=S.E
 DO N=1.KZMWD
 DO M=1.KYMWD
  DXVOL_T(M,N,I) = W1R(M,N)
...
  END DO
 END DO
END DO
!$acc loop collapse(3)
DO I=S.E
 DO N=1, KZMWD
  DO M=1.KYMWD
VOLNON T(M,N,I) = \&
-DXPHI 17(M,N,I)*DYVOL T(M,N,I)+DYPHI T(M,N,I)*DXVOL T(M,N,I) &
+DXPSFT(M,N,I)*DYCURTT(M,N,I)-DYPSITT(M,N,I)*DXCURTT(M,N,I)
  END DO
 END DO
END DO
!$acc end kernels
!$acc update host(DENSNON T,TEMENON T)
call CALC RADLOSS(30)
call CALC_SOURCE(IPTURB)
...
```





Elapse time has been reduced from 21.18 sec to <u>9.37</u> sec for 1 GPU case (~1.81 times faster than CPU case)

Elapse time has been reduced from 9.37 sec to <u>4.34</u> sec for <u>4 GPU case</u> (~<u>3.91 times</u> faster than CPU case)

Benchmark on ATOS BULLSEQUANA X2135 (NVIDIA A100)

LEONARDO Booster, CINECA

Model	Atos BullSequana X2135 "Da Vinci" single-node GPU blade
Racks	116
Nodes	3456
Processors	single socket 32 cores Intel Ice Lake CPU 1 x Intel Xeon Platinum 8358, 2.60GHz TDP 250W
Accelerators	4 x NVIDIA Ampere GPUs/node, 64GB HBM2e NVLink 3.0 (200GB/s)
Cores	32 cores/node
RAM	512 (8x64) GB DDR4 3200 MHz
Peak Performance	about 309 Pflop/s
Internal Network	DragonFly+ 200 Gbps (NVIDIA Mellanox Infiniband HDR) 2 x dual port HDR100 per node
Storage (raw capacity)	137.6 PB based on DDN ES7990X and Hard Drive Disks (Capacity Tier) 5.7 PB based on DDN ES400NVX2 and Solid State Drives (Fast Tier)

	CPU 32MPI	4MPIx4GPU (1 node)	8MPIx4GPU (1 node)
matx	0.37 (sec)	1.54	0.88
vect	5.62	23.86	12.96
rhs	11.90	6.58	5.08
fft	4.24	1.88	0.61
com	0.95	1.31	0.76
pus	27.49	164.95	83.13
gthr	4.90	10.59	6.23
com	2.90	6.24	2.55
totoal	93.17	241.11	146.96

The sparse matrix inversion part (LAPACK) is expensive.

GPU Communications between Nodes

In TMRHS (rhs) subroutine, MPI communication is performed using NVLINK+RDMA. For comparison to communication via host memory, 2 nodes are used for benchmark.



!\$acc host_data use_device(VOLNONK_T,VOLNONK)

CALL MPI_ALLTOALLW(VOLNONK_T, IRCNT, IRDISP, IRTYPE, & VOLNONK, ISCNT, ISDISP, ISTYPE, COMM1D, IERR)

Tunning by NVIDIA

Optimization I

Diagnosis

Relatively small size of 2D FFT is executed many times which implies low level of GPU utilization rate

Countermeasure

Pxovison#5G#IWv#duh#h{hfxvhg#rq#JSX#

Optimization II

Diagnosis

Pre-process and post-process of Alltoall communication are not parallelized

Countermeasure

Parallelization by OpenACC and execution on GPUs

Optimization III

Diagnosis

Data copy between CPUs and GPUs before and after Alltoall communication

Countermeasure

Direct Alltoall communication by NCCL

Test Environment

DGX H100, 4GPUs

Execution conditions

4MPI x 4GPU

Asis: cuFFT on GPU, Alltoall communication, pre-process and post-process on CPU Batched FFT: optimization I and optimization II NCCL: optimization I, optimization II and optimization III

	Asis	Batched FFT	NCCL
Time: rhs	24.05	10.69	3.44

Summary and Discussion

The optimization of R5F code for GPU version is performed using OpenACC.

- \checkmark TMRHS (rhs) subroutine is tuned which is the most expensive part of the code.
- ✓ We have reduced communications between CPU and GPU in between cuFFT calls.
- ✓ In addition, we have extended 4 GPU use in 1 node.

As the result, elapse time of subroutine TMRHS is reduced more than half.

Supercomputer LEONARDO in CINECA is also used for benchmark of R5F.

- ✓ Sparse matrix inversion part (pus) by LAPACK is most expensive part, it should be ported in GPU side.
- ✓ It is shown that NVLINK gives better performance for MPI communication compared with IB.

In future work, we should implement

fftw_plan_many_dft => cufftPlanMany ZGBSV => cuSPARSE